# 10 ABSTRACT OF THE DISCLOSURE

The grammar of this invention is an approach to regular expressions which introduces advantages to programmers who use regular expressions for scanning, searching, and tokenizing text:

5    ➢    Allowing composition of regular expressions (patterns) through the standard C/C++ operators (using C/C++ precedence and associativity), thus appealing to a huge audience of programmers already familiar with that style of grammar.

     ➢    Generalizing the negated character-class (so familiar to Perl programmers) in a way that quite dramatically allows ANY pattern composition to be used for searching (the "subjunc-

10      tive" binary composition).

     ➢    Integrating into an elegantly simple grammar form ("do-pattern") the ability to create arbitrary side-effects of tokenization, accomplished in prior art through a cumbersome combination of tokenizing expressions and parse trees, such as in the grammar style of Lex-Yacc.

     ➢    Generalizing the capture-to-variable feature (as seen in Perl), allowing the capture of por-

15      tions of the stream (match sub-expressions) into any variable current in the scope of the regular-expression.

     ➢    Allowing the parameterization of production rules, as templates, which allow similar (in form) regular expressions to be written as multiple instantiations of the same production rule template (accomplished via *in* params).

20    ➢    Further allowing parameterization of production rule templates to extend to the capture output of the expressions (accomplished via a combination of "do-patterns", "capture-patterns", and *out* or *in/out* params).

     ➢    Creating novel support algorithms (to accomplish the above) not seen in any texts on finite automata.

The following table of contents is not intended for publication but is provided as an aid to prosecution of the present application since all page number references will change with publication.

# Table of Contents